PCT

WORLD INTELLECTUAL PRO
International

INTERNATIONAL APPLICATION PUBLISHED UNDE

WO 9603741A1

| (51) International Patent Classification 6 :<br><br>G10L 5/06, 7/08, 9/06, 9/18 | A1 | (11) International Publication Number: WO 96/03741 |
|---|---|---|
| | | (43) International Publication Date: 8 February 1996 (08.02.96) |

(72) Inventors: PFISTER, Henry, L.; 4455 Torrance Boulevard #147, Torrance, CA 90503 (US). SMITH, George, W.; Suite 200, 23842 Hawthorne Boulevard, Torrance, CA 90505 (US). TSUCHIYA, Masahiro; 1585 Kaminaka Drive, Honolulu, HI 96816 (US).

(54) Title: SYSTEM AND METHOD FOR FACILITATING SPEECH TRANSCRIPTION

(57) Abstract

The invention provides a system and method for facilitating speech transcription which accepts continuous speech from any of a variety of conventional devices capable of converting spoken words to electromagnetic signals, including microphones or telephones and, if the input signal is analog, converts the input signal from analog to digital format. The digitized signal is then processed in the time and frequency domains to extract spectral speech features which are used to match the input speech with associated phonemes. According to the invention, possible word choices may be extrapolated from the associated phonemes, and visually displayed in textual representations. The visually displayed text may then be edited and processed into final form. Systems and methods for facilitating the invention are also disclosed.

1

## DESCRIPTION

## System and Method for Facilitating Speech Transcription

## Background of the Invention

The field of the invention is speech recognition and
transcription systems. Conventional speech recognition
systems have used several independent approaches with
5   limited success. One approach models the vocal tract,
articulation, and acoustic production of speech. A second
approach models the acoustic waveform and the spectrum of
human speech and uses signal processing methods. A third
approach models the human ear and its detection mechanisms
10  through neural network methods. A fourth approach
analyzes phonetic features, perception, and linguistic
models of human speech.

There is a substantial difference in the performance
and capability of the systems which currently exist.
15  Initially, these systems accepted only discrete speech
since the then nascent technology was unable to distin-
guish between individual words in continuous speech. This
limitation required users to pause between each word. As
a result, discrete speech systems proved impractical to
20  use due to the onerous speech adjustments required.

Efforts to improve upon discrete speech systems
resulted in the development of systems which could
identify one or more designated words within a continuous
stream of words. Although these systems could not
25  transcribe continuous speech, they were adequate for
applications with limited functionality, such as remote
control programming for audio-visual equipment.

A still higher level of complexity achieved by speech
recognition systems was the recognition of every word
30  within a sentence. However, these systems required that
sentences be constrained within certain grammatical
boundaries.

2

Current efforts in the field of speech recognition
have focused on developing speech processors which accept
continuous speech while attempting to overcome the above
limitations and constraints.  These systems perform inade-
5  quately when confronted with identifying word boundaries,
providing correct spelling, and resolving other word
ambiguities.

For example, word boundary issues can arise when any
subset of sequential syllables within a multi-syllable
10  word is itself a word.  Word boundary problems also arise
when the ending syllable or syllables of one word can be
joined with subsequent syllables to form another word.
Homonyms and multiple spellings of one word also give rise
to ambiguities.  Therefore, fully automated, accurate end-
15  to-end continuous speech recognition involves intensive
data analyses, speech understanding and logical reasoning
capabilities.   Currently available continuous speech
processors thus suffer from transcription errors and other
limitations as the requisite technology is currently
20  prohibitively expensive and uneconomical for commercial
purposes.

Summary of the Invention

The system and method according to the present
invention facilitates speech transcription of normally
25  spoken continuous speech without sacrificing accuracy in
transcription.  Speech is digitally processed to extract
its spectral features.  These features are then used to
distinguish individual phonemes and to generate a string
of equivalent machine-readable phonetic symbols.   The
30  string of machine-readable phonetic symbols is processed
to identify possible word boundaries for each spoken word
within the string.  In one embodiment, this process is
performed one spoken word at a time.  In this embodiment,
the possible word choices for each spoken word, as
35  delineated by the identified word boundaries, are visually
displayed.  The term "word choice" as used herein through

3

this application represents both singular and multiple word choices as both possibilities may exist and option- ally may be displayed for each spoken word. The represen- tation format can vary depending on the particular
5   language used. For example, if the speech is in the Japanese language, word choices may be represented as a string of Hiragana and Katakana symbols, each correspond- ing to a machine-readable phonetic symbol, or as Kanji characters, or as a combination of Kanji, Hiragana and
10  Katakana letters. Word choices in English, German, French, or other languages may be represented as alpha- numeric text or other appropriate representations. In each case, the proper word choice corresponding to each spoken word can be readily selected, thereby ensuring
15  accurate transcription. As the technology develops, it would also be within the scope of the invention to repre- sent the possible word choices in different languages than that dictated. Moreover, the possible word choices can be displayed in various orders, for example, by alphabetical
20  order, by order of probability, or by order of syllabic length, thereby facilitating word selection. Other display orders, including those discussed later, would be apparent to those skilled in the art and are within the scope of the invention.
25      After the first spoken word within the machine- readable phonetic symbol string is resolved in the above manner, the system proceeds to the next word within the machine-readable phonetic symbol string and iterates the aforementioned process. The process may also use adaptive
30  feedback after each selected word choice to automatically determine the starting word boundary for the next word within the string. Thus, the process always starts at the beginning of the next word, thereby further increasing accurate determination of word boundaries. Moreover,
35  editing tasks such as punctuation, margins, paragraphs, and capitalization, can be performed during this process. Other embodiments process the input speech in larger

4

blocks of speech.  For instance, an embodiment can process
the input speech for any block of speech which is brack-
eted by silence intervals, such as a phrase or sentence at
one time.  Word choices are then presented for each word
5    within the block of speech at one time.

In contrast to conventional methods of speech
recognition which attempt to resolve ambiguities, the
invention recognizes that the limitations of currently
available conventional technology render this task
10   commercially impractical.  Accordingly, in one embodiment,
the system according to the invention truncates the
conventional process and presents possible solutions to
the ambiguities in an environment wherein the correct
alternative may be readily selected.  Thus, according to
15   the invention, ambiguities may be readily and accurately
resolved while repetitive tasks, i.e. collecting, convert-
ing and analyzing massive speech data and searching for
words in a vocabulary set are performed automatically.

Accordingly, it is an object of the present invention
20   to provide a system and method for facilitating continuous
speech transcription.  The invention has the further
advantages of being both language and speaker independent,
in that the invention accommodates any speaker and any
spoken language which can be transcribed.  Other and
25   further objects and advantages will appear hereinafter.


Brief Description of the Drawings

FIG. 1 is a logical diagram of a preferred embodiment
of the invention.

FIG. 2 is a block diagram of speech acceptance and
30   analog-digital hardware suitable for practicing the
invention.

FIG. 3 is a block diagram of a spectral processor
suitable for practicing the invention.

FIG. 4 is a block diagram of a phoneme labeller
35   suitable for practicing the invention.

5

FIG. 5 is a block diagram of a speech re-synthesizer suitable for practicing the invention.

FIG. 6 is a graph of the performance of a pre-emphasis filter which is suitable for practicing the

5   invention.

FIG. 7 is a graph and equations showing computation of acoustic power.

FIG. 8 is a graph and equation showing peak-to-peak pitch estimation.

10  FIGs. 9A-9B are graphs and equations showing frequency domain processing.

FIGs. 10A-10B are an example of forward pass labelling.

FIG. 11 is an example of backward labelling.

15  Appendix A describes spectral distortion techniques.

Appendix B describes Mel-scale filters.

Appendix C describes weighting techniques.

Appendix D describes forward pass and backward labelling.

20  Appendix E describes ranking word candidates.


Detailed Description of a Preferred Embodiment of the Invention

While the invention will be described in terms of preferred embodiments, it will be appreciated that other

25  embodiments are possible within the scope of the invention.

As shown in FIG. 1, a preferred embodiment of the system according to the invention has a language, user, and mode selector 2, analog-to-digital convertor 4 inter-

30  acting with an input device 6, a spectral processor 8, phoneme labeller 10 interfacing with phoneme models 14 and preferably, training data 16, a software pre-processor 11, and a word processor 12. In another preferred embodiment, the system may also include a speech re-synthesizer 13.

35  According to the method of the invention, the speaker dictates speech into the system through the input device

6

6, which may be a microphone, telephone, other line input, or any input device which converts speech into electromagnetic signals. The term, 'electromagnetic signal', as used throughout this application, includes any form of

5  electrical signals, including microwaves, radio signals, and television signals, as well as optical and other signals such as infrared and laser.

At start-up, the speaker identifies him or herself and selects a language and an operational mode. Three of

10  the possible operational modes are training, dictation and display and editing. These three modes are those known to the inventors at the present time. However, other operational modes may become apparent to one skilled in the art as the technology develops and would also be within the

15  scope of the invention. As explained in detail later, the dictation and display and editing modes may be operated concurrently.

In the training mode, the speaker dictates speech into the system. The system then analyzes the speech to

20  generate data representing the speaker's distinctive voice characteristics. This data is stored as training data and is later used by the system to identify the speaker's words when the system is operated in dictation mode. While not required to practice the broad scope of the

25  invention, completion of the training mode in a preferred embodiment greatly facilitates the efficient practice of the invention.

In dictation mode, the speaker dictates speech into the system, using the input device 6. In one embodiment

30  of the invention, the speaker may dictate the speech at a prior time using conventional means such as a tape recorder and the audio or electromagnetic output of the conventional means can be transmitted to the input device 6. The input speech is passed through the analog-digital

35  converter 4, then through the spectral processor 8, and then through the phoneme labeller 10. The output of the phoneme labeller 10 can either be stored for later

retrieval and processing, or can be sent to the software pre-processor 11. It will be appreciated that elements of the system may be combined or rendered unnecessary as technological capabilities improve, while remaining within 5 the scope of the invention.

In the display and editing mode, editing and other functions may be performed through the software pre-processor 11 and word processor 12, using keyboard, mouse, digitizing pen, voice commands, or any other input-10 positioning device activated contemporaneously or at a later time by the speaker, by a subsequent user, or by a suitable software/computer combination. The software pre-processor 11 receives the output of the phoneme labeller 10 and works in conjunction with the word processor 12 to 15 display possible word candidates for each spoken word. The displayed word choice is selected by the speaker, a subsequent user, or by an appropriate software/computer combination. The speaker or a subsequent user may also manually input an alternative word choice. This process 20 can then be repeated for the next and subsequent words within the input speech, until the entire input speech has been transcribed. In one preferred embodiment, if a user cannot understand any segment of the spoken input, the system can include a speech re-synthesizer 13 which can 25 reproduce the spoken words. It will be appreciated that within the scope of the invention, the speech input may be stored for later processing. In this embodiment, the speech re-synthesizer can facilitate accurate transcription.

30 Preferred embodiments of the invention are now described in detail.

Training Mode

Preliminary training is preferably performed to initialize the system for each particular speaker or for 35 each time that a particular speaker uses a new language. During this initialization period, the speaker-specific

8

training data 16 is established to tune the system to the
speaker's individual speech characteristics for the
particular chosen language. In this preferred embodiment,
the system thus adapts to and accommodates a diverse set

5    of speakers regardless of gender, age, accent, dialect,
language or any other factor which can contribute to a
difference in the pronunciation of any particular word.
Initialization only needs to be performed once for any
particular speaker for each particular language.

10       The speaker-specific training data 16 can be
established independent of the invention or a system
according to the invention may be used. In embodiments
where the speaker generates training data by using the
system, the speaker dictates a pre-established reference

15   set of words into the system. Various reference sets of
words can be used for any particular language. Methods
for determining these reference sets are well-known to
those skilled in the art and are also set forth in the
publication, Fundamentals of Speech Processing by Rabiner

20   and Juang, Prentice Hall, 1993, which is incorporated
herein by reference. The dictated speech is then pro-
cessed through an analog-to-digital convertor 4, spectral
processor 8, and phoneme labeller 10 shown in FIG. 1. The
analog-to-digital convertor 4 and spectral processor 8

25   process the dictated speech in the same manner during a
training mode as during a dictation mode and are described
in more detail later for the dictation mode.

The phoneme labeller 10 is then used to generate
training data 16 by using a spectral distortion measure to

30   compare the speaker's speech against reference speech
samples of the pre-determined set of words. Spectral
distortion techniques are well-known to those skilled in
the art and are also documented in Discrete-Time
Processing of Speech Signals by Deller, Proakis, and

35   Hansen, which is incorporated herein by reference and is
also attached as Appendix A. Training data can also be

generated independent of the invention through other techniques well known to those skilled in the art.

Preferably, the training data 16 is stored into a data base. The data base may be organized in a variety of
5 suitable structures; however, a preferred embodiment uses the same structure as in DARPA TIMIT Acoustic-Phonetic Continuous Speech Corpus CD-ROM NISTIR 4930, NTIS, February, 1993, which structure is shown in Example 1:

Example 1:
10 database hierarchy:
/<corpus>/<usage>/<dialect>/<sex>/<speaker>/<word id>.file
corpus := language
usage := train/test
dialect := $d_1 \ldots d_n$
15 sex := m/f
speaker := (initials)(digit)
word id := <text type><word number>
text type := sa/si/sx
word number := 1...n words
20 file := wav/txt/wrd/phn

Dictation Mode
Speech Acceptance and Analog-Digital conversion.

In ordinary use, that is, after the system has received and processed training data for the speaker and
25 the chosen language, the system of this preferred embodiment is used in dictation mode. As indicated in FIG. 2, speech is received into the system via any of a variety of conventional devices, including but not limited to static or dynamic microphones 18, telephones 22, or other line
30 input 20. As will be appreciated by those skilled in the art, speech input could also be transmitted and received by other technology via electromagnetic signals.

If microphones are used, they preferably have noise limiting capability. As shown in FIG. 2, appropriate
35 amplifiers 24 and interfaces 26 are used to connect the

10

receiving devices with a low-pass analog filter 28 and to ensure a fixed level output compatible with digital equipment. These devices are readily available. Moreover, one skilled in the art would know how to choose appropriate
5    amplifiers and interfaces. The text, Digital Signal Analysis, by Stearns, Hayden 1975 can also be referenced to determine the compatibility of these devices.

Once received, the signal representing the input speech is passed through a low pass analog filter 28 which
10   eliminates frequencies that can alias into the acoustic signal during subsequent processing. One preferred embodiment uses a 3rd order Butterworth analog filter with 10 kHz and 6 dB/octave as the low-pass filter of choice.

The signal is then passed through an analog-to-
15   digital converter 30. While a variety of analog-to-digital converters may be used for this purpose, a preferred embodiment uses a Motorola 56 ADC 16 with sixteen bit resolution at 20 Khz and a continuously adjustable gain over a 20 db range.

20   Phonetic Feature Detection

The digitized signal is passed through a spectral processor, as shown in FIG. 3, to extract spectral speech features. Before describing this process in detail, it is useful to discuss the problem of environmental or
25   "background" noise which is commonly encountered in speech recognition. Two classes of problems contribute to background noise. One class is attributed to the speaker and consists of sound artifacts such as lip smacks, heavy breathing, mouth clicks, and nasal pops. Although such
30   artifacts are generated inadvertently, they often have an energy level comparable to speech. As will be further discussed, the invention preferably models speaker artifacts and detects them along with phoneme recognition so that they may be removed during subsequent processing.
35   A second class of noise problems arises from the ambient noise environment. Occasionally, non-speaker

11

generated background noise, such as a bell, whistle or other sound, may interfere with and mask the speech. To alleviate this problem, the system can use a noise limiting microphone for receiving the dictated speech.
5    Other means of eliminating background noise will be apparent to those skilled in the art.

Additionally, or independently, a running average of the ambient background noise during periods of silence can be kept so that sudden non-vocal tract background noises
10   can be easily classified and eliminated. As will be explained in more detail, a preferred embodiment uses both a noise limiting microphone and continuous modelling of background noise. These techniques avoid the problems commonly experienced by other speech processing systems
15   that pre-process the signal to remove noise, as the latter lose or distort the speech.

In a preferred embodiment, the digitized signal is first passed through a pre-emphasis filter 32 which works in connection with a noise model 34 to eliminate back-
20   ground noise. The noise model 34 is initially set to a default noise level and is adaptively updated during subsequent processing of the input speech. This adaptive feature is later discussed in connection with frequency domain processing.
25    In one preferred embodiment, the pre-emphasis filter 32 is a Finite Impulse Response (FIR) bandpass digital filter with zero phase shift and unity gain. The derivation of the bandpass coefficients for such a filter is shown in Example 2.

30   Example 2.
A 20 khz sampling rate with a bandpass from 100 hz to 6000 hz is assumed. The sample signal has a 60 hz hum component, a 7000 hz whistle, random noise, and an acoustic signal at 1000 hz.
35   Filter design characteristics:
Lower cutoff frequency: wlo = 100

12

Upper cutoff frequency:   whi = 6000

Sample period:            dt = .00005

rate:                     = 1/dt = 2 X $10^4$

Number of coefficients:   nc = 45

5  Filter coefficients computation:

whir = whi * 2 * $\pi$ = 3.77 X $10^4$

wlor = wlo * 2 * $\pi$ = 628.319

$c_0$ = dt * (whir - wlor) / $\pi$ = 1.8574

i    =    1    ..    nc;    $c_i$=    (sin(i*whir*dt)    -

10      sin(i*wlor*dt))/($\pi$*i)

Sample signal:

nn = 511   nt = 2     noise = 0.5

$f_0$ = 60    $f_i$ = 1000   $f_2$= 7000

$a_0$ = 0.5    $a_1$= 1.0     $a_2$= 0.5

15      i = 0 .. nn

$$input_i = \sum_{j=0}^{nt} a_j * \sin(i*dt*2.0*\pi*f_j) + rnd(noise) - \frac{noise}{2}$$

$output_i$ = 0

$error_i$ = 0

$signal_i$ = sin(i*dt*2.0*$\pi$*f1)

k= nc..nn-nc

$$output_k = c_0 * input_k + \left[ \sum_{j=1}^{nc} [(c_j * input_{k+j}) + (c_j * input_{k-j})] \right]$$

$$error_k = |output_k - signal_k|$$

20      The performance of the pre-emphasis filter 32 is
shown in FIG. 6.  In FIG. 6, the upper trace (a) is the
acoustic input signal.  The middle trace (b) is the result
of the pre-emphasis filter.  The third trace (c) is the

13

desired signal that is embedded in the signal. The lower trace (d) is the absolute difference of the acoustic signal and the filter output.

The signal is then sampled into blocks of data as
5   indicated by the sampler 36. Since the human vocal mechanism modulates the slowly changing speech signal onto a higher frequency sound wave, it would be necessary to sample the acoustic wave form at over 10,000 times per second to capture this speech wave. In contrast, the
10  movement of the tongue, jaw, lips and other vocal articulators change at the far slower rate of less than 100 times per second. This physical situation is exploited by grouping acoustic data into approximately one hundredth of a second blocks to isolate phonetic features and to
15  identify noise sources.

Taking the above considerations into account, a preferred embodiment samples the signal into blocks of 512 samples. Oversampling can be used to increase accuracy during subsequent processing. Thus, a preferred embodi-
20  ment uses a 25%, or 128-sample overlap between adjacent blocks. As can be readily appreciated by one skilled in the art, numerous other data block sizes and overlap ranges can be used. For example, blocks of 256, 1024 or 2048 samples may be used. Moreover, any overlap range
25  between 0 and 50% may be used between adjacent blocks.

While sampling may be accomplished with hardware, microcode, or other methods known to those skilled in the art, a preferred embodiment uses microcode. As technology advances, other sampling methods and devices suitable for
30  practicing the invention may become available.

The sampled signal is then processed in the time domain, as indicated by element 38, to extract the spectral speech features of acoustic power and the peak-to-peak pitch.

35  Acoustic power indicates the presence of speech or silence intervals that naturally occur at the end of sentences, phrases, or words. To determine acoustic

14

power, a short time energy estimate is made for each
sample block by averaging either the magnitude or square
of the signal within the block, otherwise referred to as
absolute power and squared power, respectively. Next, the
5    peak-to-peak pitch is estimated by summing the signal zero
crossings.  Although this process is well-known to one
skilled in the art, an example of time domain processing
is given in Figures 7-8.

In FIG. 7, a sample block of 512 samples is shown
10   with a mean value of $3.298 * 10^4$. Plot (a) represents the
sample block signal.  Equations (b) and (c) are used to
compute the absolute power and squared power,
respectively.

In FIG. 8, the peak-to-peak pitch is estimated by
15   summing the signal zero crossings of the estimated signal
shown in plot (d), as shown by equation (e).

The signal is also passed through acoustic band pass
filters 40 to identify the relative signal power in each
acoustic band, otherwise referred to as the spectral
20   pattern.  This information can be used later to assist in
phoneme identification.  In a preferred embodiment, the
band pass filters are selected based on the Mel scale.
The Mel scale is a logarithmic-based scale which is more
fully described in the publication <u>Advanced Algorithms and</u>
25   <u>Architectures for Speech Understanding</u>, ESPRIT Research
Report Project 26, Vol. 1, Pirani, 1990, and is herein
incorporated by reference, and is also attached as
Appendix B.

The estimated formant frequencies are determined by
30   processing the signal in the frequency domain as shown by
element 42.  An example of frequency domain processing is
given in Figures 9A-9B.  In Figure 9A, plot (a) represent
the same sample block as plot (a) in Figure 7.  The sample
block can be first tapered using a tapering function.
35   Numerous tapering functions, including a Hanning Window as
shown in Equation (b), can be used to taper the 512
elements of the sample block, and are well-known to those

skilled in the art. Tapering reduces high frequency noise in the subsequent Fast Fourier Transform (FFT). Since the signal was previously oversampled by 256 elements (128 elements at each end), tapering promotes accuracy, rather

5   than loss of data. Plot (c) represents the tapered sample block. As would be appreciated by one skilled in the art, numerous other tapering functions can be used.

An FFT is then performed on the tapered sample block and the power spectrum is estimated using the square of

10  the resulting complex frequency coefficients, as shown by equation (d) in FIG 9B. Plot (e) provides a graph of the resultant power spectrum. The log of each element is then computed, as shown in plot (f). Since human speech usually does not exist at frequencies greater than 8000

15  hz, the frequency spectrum in plot (f) can be truncated to reduce high frequency noise. The cut-off range can be as low as 7500 hz or as high as 8500 hz. In the embodiment shown in plot (g), the frequency is truncated at 7500 hz. Next, an inverse FFT is performed to determine the cep-

20  stral coefficients, as shown in plot (h). The resonant or formant frequencies of the vocal tract are then extracted from the low order peaks of the cepstrum in plot (h).

Note that if a silence interval is detected during time domain processing, the data from frequency domain

25  processing necessarily represents only ambient background noise. The system and method according to the invention update the noise model 34 with this data. Thus, the noise model 34 reflects a semi-continuous model of background ambient noise. This adaptive feedback enhances subsequent

30  processing of speech through the pre-emphasis filter 32. Thus, in this embodiment, the spectral processor extracts the following spectral speech features: acoustic power, spectral pattern, and formant frequencies.


Phoneme Identification.

35      The extracted spectral speech features are passed to a phoneme labeller 10 as shown in FIG. 4. As would be

known to one skilled in the art, the previously extracted spectral feature of relative signal power can be used to classify phonemes into one of the following phonetic categories: vowels, diphthongs, liquid or glide semi-

5  vowels, nasal consonants, voiced or unvoiced fricatives, africates, whisper and voiced or unvoiced stops. Table 1 describes a possible phoneme classification scheme for forty standard phonemes used in Western/European languages:

10  Table 1:

vowels divided into:

front -- IY, IH, EH, AE

mid -- AA, ER, AH, AX, AO

back -- UW, UH, OW

15  dipthongs - AY, OY, AW, EY

semivowels, including

liquids -- WW, LL

glides -- RR, YY

nasal consonants -- MM, NN, NG

20  stops

voiced -- BB, DD, GG

unvoiced -- PP, TT, KK

fricatives

voiced -- VV, TH, ZZ, ZH

25  unvoiced -- FF, TH, SS, SH

africates -- JH, CH

whisper -- HH

By classifying phonemes into phonetic categories, the individual phoneme identification problem can be greatly

30  simplified. For instance, under the classification scheme given in Table 1, each phonetic category contains at most five elements out of a possible forty. Moreover, these elements are usually well separated by physical features of the vocal tract so that subsequent identification of an

35  individual phoneme within a phonetic category is enhanced.

17

To identify individual phonemes for each signal block, a 'forward pass' 44 is made through each signal block and phoneme candidates are identified by comparing against existing features in the phoneme models 14. The

5   phoneme models 14 are pre-determined reference sets of data for each language. Techniques for generating these sets are well-known. For instance, the TIMIT Acoustic-Phonetic Continuous Speech Corpus CD-ROM NISTIR 4930, NTIS, February, 1993 is one such set which is readily

10  available to one skilled in the art. A score is kept for each phoneme candidate based on its probability of being the correct phoneme spoken.

In a preferred embodiment, scoring is accomplished by using the following fuzzy logic LR membership functions:

15          Left(z)  =  1/(1 + z*z)          Left side function
            Right(z) =  1/(1 + z*z)          Right side function
The function is described as a four point trapezoid (vw) and disjunction:

$$ww(a,b,c,d,y,x) = if\left(x<b, left\left(\frac{b-x}{b-a}\right), if\left(x>c, right\left(\frac{x-c}{d-c}\right), y\right)\right)$$

        mx(a,b) = if(a>b,a,b)        mn(a,b) = if(a<b,a,b)
20      w x ( a 1 , b 1 , c 1 , d 1 , a 2 , b 2 , c 2 , d 2 , y , x )      =
        ww(mx(a1,a2),mx(b1,b2),mn(c1,c2),mn(d1,d2),y,x)


These functions are well-known to those skilled in the art and are also discussed in Fuzzy Set Theory and Its Applications by Zimmermann, 2d, Kluwer, 1991, which is

25  herein incorporated by reference. In this preferred embodiment, each phoneme candidate feature and sample feature are used to form a normalized membership function of the sample in the phoneme data base. Thus, all phonemes have a possibility of assignment in each sample

30  block. The features can be weighted to maximize the probability of phoneme identification. The weighting

18

process can be accomplished through a variety of tech-
niques which are well-known to those skilled in the art,
some of which are described in <u>Discrete-Time Processing of
Speech Signals</u>, which are incorporated herein by reference
5   and are also attached as Appendix C.   Preferably, the
determination of the weighting function is based on
experimental scoring derived from the training data 16 for
the speaker.   As discussed later, the weighting function
can be updated through adaptive feedback from the user.
10  An example of forward pass labelling is given in FIGs.
10A-10B.

As previously mentioned, speaker artifacts and
sudden, loud background noises can interfere with the
input speech so that the system is unable to find a
15  reasonable phoneme match for the received sound.   In a
preferred embodiment, such interference is identified and
the speaker is requested to repeat the masked speech.

When the extracted spectral speech features indicate
a silence interval, a "speech segment" bracketed by the
20  silence interval and its initial starting point has been
defined.   This event triggers backward labelling 46 of the
most likely phoneme candidates for the speech segment.
Simultaneously, spectral processing 8 and forward pass
labelling 44 continue for the next and subsequent speech
25  segments within the speech signal.

A silence interval could be defined as any time
interval longer than 1/10 (one-tenth of a second) in which
the ambient noise level does not exceed 10 dB over the
average ambient background noise level.   Shorter time
30  intervals and smaller decibel increases may also be used.
For instance, a shorter time period for speakers who
dictate too rapidly to pause for more than 1/10 of a
second, or for speakers who speak too softly to develop a
10 dB increase over the ambient background noise level.
35  Thus, it is also within the broad scope of the invention
for embodiments to define silence intervals as short as
1/15 of second, or to require only a 9 dB increase.   Other

19

embodiments may require over 12 dB increases for speakers who are particularly loud. Other combinations are readily apparent to those skilled in the art.

Backward labelling for these subsequent speech
5    segments is initiated once processing for the instant speech segment is complete. Notably, partitioning of the speech into speech segments bracketed by silence intervals promotes accuracy and facilitates parallel processing. Moreover, bracketing further facilitates the display of
10   full phrases in embodiments which process speech in blocks of phrases at one time. These embodiments are described later.

Backward labelling identifies the N most possible phoneme sequences within the speech segment. Each phoneme
15   candidate is ranked by its possibility of assignment in each sample block, based on conventionally known methods of phoneme identification. These methods are discussed in <u>Discrete-Time Processing of Speech Signals</u>, and are incorporated herein by reference, and are also attached as
20   Appendix D. As shown in FIG. 11, the candidates are sorted to maximize the possibility of phoneme identification.

Once the phoneme candidates are identified, the system generates an equivalent machine-readable phonetic
25   symbol for each phoneme candidate, as shown by element 48. This process utilizes look-up tables 15 which associate machine-readable phonetic symbols with each phoneme model 14. The machine-readable phonetic symbols are composed of machine-readable codes. Thus, by cross-referencing each
30   phoneme candidate with the equivalent phoneme model and its associated machine-readable phonetic symbol, a corresponding machine-readable phonetic symbol is generated for each phoneme candidate. Although various machine-readable codes can be used, a preferred embodiment uses
35   ASCII codes.

20

Output to Software Pre-Processor

        In one preferred embodiment, the machine-readable
phonetic symbols are directly output to a software pre-
processor 11 as shown in FIG. 1.   The software pre-
5   processor 11 can be separately joined to or integrated
with a word processor 12 running, for example, on a
personal computer, on stand-alone equipment or as part of
a network or central computer system.   The software pre-
processor 11 may be connected to the phoneme labeller 10
10  via cable connection or any conventionally available
wireless communication device such as infrared trans-
mittal, laser, radio, microwaves, or television.   The
process then continues with the editing mode described
later.

15      In another preferred embodiment, the machine-readable
phonetic symbols can be stored on any conventionally
available device which stores machine-readable data, such
as, for example, computer data disks, hard drives, flash,
static or dynamic memory, tape, or CD-PROM.   Once the
20  machine-readable phonetic symbols are stored, the process
can be re-started at a later time by loading the machine-
readable phonetic symbols into the software pre-processor
11 and continuing with the display and editing mode
described below.

25      Another preferred embodiment automatically stores the
machine-readable phonetic symbols while outputting the
machine-readable phonetic symbols to the software pre-
processor 11.


Display and Editing Mode

30      Once the software pre-processor 11 starts to receive
the machine-readable phonetic symbols corresponding to the
last spoken speech, the system can be operated in the
display and editing mode, utilizing the software pre-
processor 11, word processor 12, and speech re-synthesizer
35  13, as described below.   In one embodiment, the display
and editing mode occurs after the speaker has finished

21

dictating, as the machine-readable phonetic symbols which
represent the dictated speech can be stored for later
processing in the display and editing mode. However, it
is important to note that the speaker may continue dictat-
5   ing while words spoken earlier are being processed through
the display and editing mode. For example, a secretary
may operate the display and editing mode while the speaker
continues to dictate. Alternatively, the speaker may
pause during dictation to edit prior dictation.

10  <u>Phonetic Symbol String Editing</u>
        In a preferred embodiment, the software pre-processor
11 provides an option to edit the machine-readable pho-
netic symbol string. This option is particularly useful
for written languages which can be readily comprehensible
15  as phonetic symbols, such as Japanese Hiragana and
Katakana. To edit the machine-readable phonetic symbol
string, the software pre-processor 11 outputs the machine-
readable phonetic symbol string to the word processor 12,
where it is displayed.
20      If a wrong machine-readable phonetic symbol or set of
symbols is detected by the speaker or subsequent user, the
symbol or symbols can be manually overridden by keyboard,
mouse, digitizing pen, voice command or any other input-
positioning device such as a track ball, joystick or
25  touchscreen. For voice-input, the user manually selects
the subject symbol or symbols using either a keyboard,
mouse, digitizing pen, or other input-positioning device,
and then re-dictates the input speech as desired. This
speech is then processed as before, i.e. through the
30  analog-digital convertor 4, spectral processor 8, phoneme
labeller 10 and software pre-processor 11 to update the
machine-readable phonetic symbol string.
        If the input speech generates more than one possible
machine-readable phonetic symbol match during backward
35  labelling 46, each possibility can be displayed so that
the speaker or subsequent user can readily select the

22

correct machine-readable phonetic symbol. A preferred embodiment only displays the two most likely machine-readable phonetic symbols.

     Thus, the machine-readable phonetic symbol string may
5 be edited using the software pre-processor 11 and word processor 12. Certain languages, such as English, German, and French, differ from Japanese in that no corresponding syllabic symbols, such as Hiragana and Katakana, exist for the former. In these cases, a universal phonetic symbol
10 representation can be used. Alternatively, the speaker or subsequent user can opt to skip machine-readable phonetic symbol string editing altogether and proceed directly to locating word boundaries.

## Identifying word boundaries

15     The software pre-processor 11 identifies word boundaries within the machine-readable phonetic symbol string. Each substring of machine-readable phonetic symbols delineated by the identified word boundaries is combined into a possible word choice, which is represented
20 in· machine-readable code. While this process can be performed for the entire machine-readable phonetic symbol string or any segment thereof at one time, a preferred embodiment identifies word boundaries on a word-by-word basis. In this preferred embodiment, the software pre-
25 processor 11 starts at the beginning of the machine-readable phonetic symbol string and identifies possible word boundaries within the string <u>for the first spoken word</u>. Identification of possible word boundaries for subsequent spoken words occurs later, and can be
30 accomplished, in part, by using adaptive feedback to determine the next starting word boundary. This embodiment is described in detail later.

    As an example of identifying word boundaries for Japanese speech, the software pre-processor 11 identifies
35 every Kanji character which can be represented by any substring of machine-readable phonetic symbols with a

23

beginning word boundary at the start of the machine-readable phonetic symbol string. These Kanji characters are each considered a possible word choice. Techniques for combining machine-readable phonetic symbols into words

5  are commonly known. For instance, commercially available Japanese word processors operate on these principles to convert Hiragana and Katakana symbols into Kanji characters.

The possible word choices delineated by the

10 identified word boundaries are then ranked in order of probability by reference to linguistic usage, including such factors as grammatical and contextual syntax. Techniques for performing this task are well-known to those skilled in the art and are also described in the

15 aforementioned text, <u>Discrete-Time Processing of Speech Signals</u>, which are herein incorporated by reference and are also attached as Appendix E. Alternatively, the possible word choices may be ranked by alphabetical or reverse alphabetical order, or by increasing or decreasing

20 syllable length, where syllable length is measured by the number of syllables within a word.

An alternative method of ranking the word choices is with reference to the prior usage of words in the dictated speech. In this method, previously dictated and selected

25 words, as later described, are accessed and if matches are found, the word choices presented reflect the previous usage. The word choices can then be ranked, for example, by their most recent usage, by the time elapsed between the current words and their prior usage (i.e. age), or by

30 frequency of usage.

The possible word choices are then output, in machine-readable form, to the word processor 12, or are stored for future use.

<u>Visually Representing Possible Word Choices</u>

35  The possible word choices are then visually represented. In a preferred embodiment, the word processor 12

possesses all the features commonly available on commer-cial word processors. The word processor 12 displays the possible word choices in representations readily compre-hensible to the speaker or a subsequent user. The
5  representation format may vary depending on the language used. For instance, in Japanese, the word processor can represent each possible word choice as a Kanji character, Hiragana and Katakana representations, or a combination of the three. For English, German, French and other similar
10  languages, the possible word choices can be represented as alphanumeric text or other appropriate representations.

Numerous display orders are also possible. For instance, the possible word choices can be displayed in order of probability, alphabetical order, reverse alpha-
15  betical order, by increasing or decreasing syllable length, where syllable length is measured by the number of syllables in a word, by the most recent usage, by the time elapsed between the current word and the word choice's prior usage (i.e. age), or by frequency of usage.

20  In a preferred embodiment, the most probable word choice is displayed separate from the remaining word choices and the remaining possible word choices are then displayed as two independent sequences. One of these sequences ranks the remaining word choices by increasing
25  syllable length and the other sequence is ranked by alphabetical order.

There are numerous methods by which the most probable word choice can be displayed separate from the remaining words. For instance, the most probable word choice can be
30  displayed in the first line and the two independent sequences of remaining words as separate columns below the first line. Alternatively, the most probable word choice can be displayed in **bold** format as compared to the remain-ing words. Other variations are readily apparent to one
35  skilled in the art.

25

## Word Selection

The word processor 12 further provides the speaker, a subsequent user, or a suitable software/computer combination with functionality to readily select any displayed
5    word choice or to manually input a word.  These options may be similar to spell checking features commonly available on conventional word processing software packages. The selection process may be performed by any one of numerous means, including but not limited to voice com-
10   mand, mouse, digitizing pen, keyboard or any other input-positioning device.  Insertion of punctuation and similar tasks can be accomplished at this time either manually, or through macros or voice commands.

A preferred embodiment further provides an option
15   whereby the speaker or a subsequent user may call up the machine-readable phonetic symbol substring associated with any displayed word.  The string may then be edited in the same manner as before, after which the software pre-processor 11 re-translates the modified machine-readable
20   phonetic symbol substring into a word and displays the word for further editing.  Alternatively, the user may manually select word boundaries within the recalled machine-readable phonetic symbol substring, which the software pre-processor 11 then uses to re-combine the
25   machine-readable phonetic symbols into words which are displayed for acceptance or further editing.

It is recognized that someone other than the original speaker (referred to herein as "subsequent user") may use the system to actually perform the editing process.  Thus,
30   a preferred embodiment of the system has the additional capability of re-synthesizing the sound represented by any speech segment so that the user can hear what speech was actually dictated.

As shown in FIG. 5, the speech re-synthesizer first
35   uses a digital signal generator 50 which receives machine-readable phonetic symbols from the software pre-processor 11 and utilizes the look-up tables 15, and if desirable,

26

training data 16, to convert the machine-readable phonetic
symbols into a digital signal. This digital signal is fed
to a digital-to-analog convertor 52, and then to a speaker
driver 54 connected to an audio speaker 56. Each of these
5   devices is commercially available.

Accordingly, the speaker need not personally perform
the editing process since any ambiguities can be under-
stood by any subsequent user, especially in embodiments
incorporating re-synthesized speech processes.  It is
10  recognized that this particular re-synthesis feature can
be implemented for uses other than speech transcription.
For instance, storing speech as an analog signal consumes
a relatively large amount of storage media.  Even a
digitized speech signal requires a substantial amount of
15  storage space.  However, storing speech as machine-
readable phonetic symbols requires relatively little
space.  Thus, the above speech re-synthesis technique
constitutes an ideal speech compression method which is
also readily adapted to telephone answering machines and
20  other voice-storage applications.


Adaptive Feedback

In a preferred embodiment, the word processor 12
utilizes adaptive feedback so that information as to a
selected word is incorporated into the selection of the
25  next starting word boundary.  After the correct word
choice has been selected from the possibilities presented,
this information is fed back to the software pre-processor
11 so that the software pre-processor will begin at the
next machine-readable phonetic symbol.  Thus, subsequent
30  word boundary identification is enhanced since the
beginning word boundary has already been determined.

In addition to resolving word boundaries, each
selection of an ambiguous word necessarily resolves
ambiguities in phoneme identification.  A preferred
35  embodiment exploits this situation by also using adaptive
feedback so that user selection information can be used to

27

update the weighting functions for subsequent phoneme labelling. Thus, this feature allows the system and method according to the invention to further adapt to a speaker's particular speech characteristics.

5      It is recognized that this feedback capability may not always be desirable. For instance, a speaker's speech characteristics may vary from time to time, due to general health, stress or other factors. Feedback may be undesirable in these instances, since the weighting functions

10    will be updated with anomalous data. Therefore, unlike feedback for determining starting word boundaries, an option is provided whereby feedback to the weighting functions can be disabled.

In practice, feedback to the weighting functions is

15    preferably used only for a short period of time, depending on the consistency of an individual speaker's speech. After the weighting functions have been updated during this period, feedback to the weighting function is disabled. Feedback may be subsequently enabled to update the

20    weighting functions if a speaker's speech characteristics have changed over a period of time, for instance, as a speaker grows older.


## Subsequent Word Resolution

Once the first word is resolved in the above manner,

25    the process is iterated for the next and subsequent words. As previously discussed, once the first word has been selected, the software pre-processor 11 begins at the next machine-readable phonetic symbol within the machine-readable phonetic symbol string to identify possible word

30    boundaries for the next word.


## Other Embodiments for Longer Speech Blocks

While a preferred embodiment analyzes the machine-readable phonetic symbol string on a word-by-word basis, it is within the scope of the invention to analyze larger

35    blocks of speech at one time, for instance, sentences or

28

phrases, or any blocks of speech which are bracketed by silence intervals.

These embodiments use the same principles as in a preferred embodiment. For instance, if the speech is
5 processed in a phrase or sentence block, the software pre-processor 11 determines word boundaries for all the words in the entire phrase or sentence, which is then displayed through the word processor 12. The user is then given the option of validating the entire phrase or sentence as
10 correct, or manually editing any word boundary within the displayed block. In the latter case, the software pre-processor 11 re-identifies word boundaries for the remaining words within the displayed block, since the user's editing may change subsequent word boundaries
15 within the sentence. The updated block is then re-displayed for acceptance or further editing.

Accordingly, a system and method for facilitating speech transcription has been disclosed. Although pre-ferred embodiments are described to process continuous
20 speech, it will be appreciated that the scope of the invention will also accommodate discrete speech. It is further apparent that a preferred embodiment may be implemented in any spoken language which can be tran-scribed and is also speaker-independent. Many other
25 embodiments are easily recognized by one skilled in the art and are within the scope of the invention. For instance, the invention may be implemented to process speech on a phrase-by-phrase or sentence-by-sentence basis. A person skilled in the art could readily modify
30 the system to accept digital input, as from an audio stereo system, by omitting analog-to-digital conversion and making appropriate modifications to pre-process the digital signal. Although currently not commercially feasible, another embodiment can use a computer, such as
35 a Sun 10 Workstation, with sufficient computing capability to perform the editing and selecting process. The

invention, therefore, is not to be restricted except in the spirit of the appended claims.

30

<u>Claims</u>

1.      A system for facilitating speech transcription
comprising:

    (a)    a device which receives a signal representing
5  speech;

    (b)    an analog-to-digital convertor in communication
with said device for converting said signal into a digital
signal;

    (c)    a spectral processor in communication with said
10  analog-digital convertor for receiving said digital signal
and extracting spectral speech features from said digital
signal;

    (d)    a phoneme labeller in communication with said
spectral processor for receiving said spectral speech
15  features, identifying phonemes from said spectral speech
features and generating corresponding machine-readable
phonetic symbols for said phonemes;

    (e)    a software pre-processor in communication with
said phoneme labeller for receiving said machine-readable
20  phonetic symbols as input and for combining said machine-
readable phonetic symbols into words; and

    (f)    a word processor in communication with said
software pre-processor for visually displaying said words.


2.      The system as in claim 1 wherein said signal
25  representing speech is an electromagnetic analog signal.


3.      A method for facilitating the transcription of
speech comprising the steps of:

    (a)    receiving a digital signal representing human
speech;

30      (b)    processing said digital signal to extract
spectral speech features;

    (c)    identifying phonemes from said spectral speech
features;

    (d)    generating corresponding machine-readable
35  phonetic symbols for said phonemes;

(e) outputting said machine-readable phonetic symbols into a software pre-processor;

(f) identifying word boundaries within a substring of said machine-readable phonetic symbols;

(g) combining said substring into a possible word choice; and

(h) visually representing said possible word choice.

4. A method for facilitating speech transcription comprising the steps of:

(a) receiving a human voice speech input as an analog signal;

(b) converting said speech input from said analog signal to a digital signal;

(c) processing said digital signal to extract spectral speech features;

(d) identifying phonemes from said spectral speech features;

(e) generating corresponding machine-readable phonetic symbols for said phonemes;

(f) outputting said machine-readable phonetic symbols into a software pre-processor;

(g) identifying word boundaries within a substring of said achine-readable phonetic symbols;

(h) combining said substring into a possible word choice; and

(i) visually representing said possible word choice.

5. The method as in claim 4 wherein said step (d) of identifying phonemes comprises comparing said spectral speech features with a reference set of phoneme models.

6. The method as in claim 5 wherein said step (d) of identifying phonemes further comprises comparing said spectral speech features with speaker-specific training data generated prior to said step (d).

32

7.    The method as in claim 4, 5, or 6 wherein said step (i) of visually representing said possible word choice comprises displaying said possible word choice as alphanumeric text.

5        8.    The method as in claim 4, 5, or 6 wherein said step (i) of visually representing said possible word choice comprises displaying said possible word choice as every machine-readable phonetic symbol within said substring.

10       9.    The method as in claim 4, 5, or 6 wherein said step (a) of receiving a human voice speech input comprises using at least one of a static microphone, a dynamic microphone, a telephone, and line input.

         10.    The method as in claim 4, 5, or 6 wherein said
15  step (a) of receiving a human voice speech input occurs in a selected language and said steps (d) - (h) and said step (i) of visually representing said possible word choice occur in said selected language.

         11.    The method as in claim 4, 5 or 6 further
20  comprising the step of ranking said possible word choice by increasing syllable length, wherein the word choice with the fewest number of syllables is ranked first.

         12.    The method as in claim 4 further comprising the step of ranking said possible word choice by decreasing
25  syllable length, wherein the word choice with the most number of syllables is ranked first.

         13.    The method as in claim 4 further comprising the step of ranking said possible word choice by alphabetical order.

33

14. The method as in claim 4 further comprising the step of ranking said possible word choice by probability in reference to linguistic usage.

15. The method as in claim 14 wherein said step (i)
5  of visually representing said possible word choice comprises displaying the most probable word choice in reference to linguistic usage.

16. The method as in claim 4 further comprising the step of ranking said possible word choice by at least one
10 of alphabetical order, reverse alphabetical order, probability in reference to linguistic usage, increasing syllable length and decreasing syllable length.

17. The method as in claim 16 wherein said step (i) of visually representing said possible word choice com-
15 prises displaying the most probable word choice, in reference to linguistic usage, separated in the display from the other choices and displaying the remaining word choices in alphabetical order.

18. The method as in claim 16 wherein said step (i)
20 of visually representing said possible word choice comprises displaying the most probable word choice, in reference to linguistic usage, separated in the display from the other choices and displaying the remaining word choices in reverse alphabetical order.

25   19. The method as in claim 16 wherein said step (i) of visually representing said possible word choice comprises displaying the most probable word choice, in reference to linguistic usage, separated in the display from the other choices and displaying the remaining word
30 choices in increasing syllable length.

20. The method as in claim 16 wherein said step (i) of visually representing said possible word choice comprises displaying the most probable word choice, in reference to linguistic usage, separated in the display
5   from the other choices and displaying the remaining word choices in decreasing syllable length.

21. The method as in claim 16 wherein said step (i) of visually representing said possible word choice comprises displaying the most probable word choice, in
10  reference to linguistic usage, separated in the display from the other choices and displaying the remaining word choices in order of probability in reference to linguistic usage.

22. The method as in claim 4, 5 or 6 wherein said
15  step (c) of processing said digital signal comprises sampling said digital signal into blocks of samples.

23. The method as in claim 22 wherein said blocks of samples have a 25% overlap between adjacent blocks.

24. The method as in claim 22 wherein said step (c)
20  of processing said digital signal comprises Fast Fourier Transform after first tapering each block of samples.

25. The method as in claim 4, 5 or 6 wherein said step (c) of processing said digital signal comprises sampling said digital signal into blocks of samples with
25  0-50% overlap between adjacent blocks.

26. The method as in claim 4, 5 or 6 wherein said step (c) of processing said digital signal comprises sampling said digital signal into blocks from 256 to 2048 samples.

27. The method as in claim 4, 5 or 6 wherein said step (c) of processing said digital signal comprises sampling said digital signal into blocks with 512 samples.

28. The method as in claim 4, 5 or 6 further
5 comprising the step of displaying said machine-readable phonetic symbols after said step (f) of outputting said machine-readable phonetic symbols.

29. The method as in claim 28 further comprising the step of editing said outputted machine-readable phonetic
10 symbols after said step of displaying said machine-readable phonetic symbols.

30. The method as in claim 29 wherein said step of editing said machine-readable phonetic symbols comprises using at least one of a keyboard, a mouse, a digitizing
15 pen, a voice command and an input-positioning device.

31. The method as in claim 28 further comprising the step of editing said word boundaries after said step of displaying said machine-readable phonetic symbols.

32. The method as in claim 4, 5 or 6 wherein said
20 step (i) of visually representing said possible word choice comprises visually representing all the words in any block of speech bracketed by silence intervals at one time.

33. The method as in claim 32 wherein a silence
25 interval is any period of longer than 1/10 of a second in which the ambient noise level does not increase more than 10 dB over the ambient background noise level.

34. The method as in claim 32 wherein a silence interval is any period of longer than 1/15 of a second in

36

which the ambient noise level does not increase more than
9 dB over the average ambient background noise level.

35. The method as in claim 32 wherein a silence
interval is any period between 1/15 of a second and 1/10
5   of a second in which the ambient noise level does not
increase more than about 10 dB over the average ambient
background noise level.

36. The method as in claim 4, 5 or 6 further
comprising the step of selecting said represented possible
10  word choice.

37. The method as in claim 36 wherein said step of
selecting one of said represented possible word choice
comprises using at least one of a mouse, keyboard, digi-
tizing pen, voice command and an input-positioning device.

15      38. The method as in claim 36 wherein said step of
identifying word boundaries comprises using adaptive
feedback from the step of selecting one of said repre-
sented possible word choice as performed on the
immediately preceding word.

20      39. The method as in claim 36 wherein said step of
selecting one of said represented possible word choice
comprises employing a suitable software/computer combina-
tion to make such selection.

40. The method as in claim 36 further comprising
25  audio reproduction of said speech input from said machine-
readable phonetic symbols.

41. The method as in claim 6 further comprising the
step of selecting said represented possible word choice,
wherein said speaker-specific training data is updated via

37

adaptive feedback from said step of selecting said represented possible word choice.

42. The method as in claim 4, 5 or 6 further comprising the step of selecting one of a previously represented possible word choice wherein at least one of said steps (a) through (i) occurs simultaneous with said step of selecting one of a previously represented possible word choice.

43. The method as in claim 4, 5 or 6 further comprising the step of storing said machine-readable phonetic symbols.

44. A method for facilitating speech transcription comprising the steps of:

(a) receiving a human voice speech input as an analog signal;

(b) converting said speech input from said analog signal to a digital signal;

(c) processing said digital signal to extract spectral speech features;

(d) identifying phonemes by comparing said spectral speech features with a reference set of phoneme models and speaker-specific training data;

(e) generating corresponding machine-readable phonetic symbols for said phonemes;

(f) outputting said machine-readable phonetic symbols into a software pre-processor;

(g) displaying said machine-readable phonetic symbols;

(h) editing said machine-readable phonetic symbols;

(i) identifying word boundaries within a substring of said achine-readable phonetic symbols;

(j) editing said word boundaries;

(k) combining said substring into a possible word choice;

38

(l)    ranking said possible word choice; and

(m)    visually representing said possible word choice.


45.    The method as in claim 44 further comprising the step (n) of selecting one of said represented possible
5    word choices.


46.    The method as in claim 45 wherein said steps (a) through (n) are repeated for a word within said human voice speech input and said step (i) of identifying word boundaries includes the use of adaptive feedback from the
10    step (n) of selecting one of said represented possible word choice which occurred for the word preceding said word within said human voice speech input.


47.    A method for storing speech as machine-readable phonetic symbols comprising the steps of:
15        (a)    receiving a human voice speech input as an analog signal;

(b)    converting said speech input from said analog signal to a digital signal;

(c)    processing said digital signal to extract
20    spectral speech features;

(d)    identifying phonemes from said spectral speech features; and

(e)    generating corresponding machine-readable phonetic symbols for said phonemes.


25        48.    A system for storing speech as machine-readable phonetic symbols comprising:

(a)    a device which receives an analog signal representing human speech;

(b)    an analog-to-digital convertor in communication
30    with said device for converting said analog signal into a digital signal;

(c)    a spectral processor in communication with said analog-to-digital convertor for receiving said digital

39

signal and extracting spectral speech features from said digital signal; and

(d) a phoneme labeller in communication with said spectral processor for receiving said spectral speech
5 features and for identifying phonemes from said spectral speech features.

49. The system as in claim 48 wherein said analog signal representing human speech is an electromagnetic signal.

10 50. A system for storing speech as machine-readable phonetic symbols comprising:

(a) a device which receives a digital signal representing human speech;

(b) a spectral processor in communication with said
15 device for receiving said digital signal and extracting spectral speech features from said digital signal;

(c) a phoneme labeller in communication with said spectral processor for receiving said spectral speech features, for identifying phonemes from said spectral
20 speech features; and for generating corresponding machine-readable phonetic symbols for said phonemes; and

(d) a storage device for storing said machine-readable phonetic symbols.

51. The system as in claim 50 wherein said digital
25 signal representing speech is an electromagnetic signal.

52. A method of speech reproduction comprising re-synthesis of stored machine-readable phonetic symbols which represent a compressed form of speech.

53. The method as in claim 4 further comprising the
30 step of ranking said possible word choice in order of the most recent usage, wherein the word choice with the least

number of words since the last usage of said word choice is ranked first.

54. The method as in claim 53 wherein said step (i) of visually representing said possible word choice comprises displaying said possible word choice in order of the most recent usage, wherein the word choice with the least number of words since the last usage of said word choice is ranked first.

55. The method as in claim 4 further comprising the step of ranking said possible word choice in order of frequency used, wherein the word previously used most often is ranked first.

56. The method as in claim 55 wherein said step (i) of visually representing said possible word choice comprises displaying said possible word choice in order of frequency used, wherein the word previously used most often is ranked first.

57. The method as in claim 4 further comprising the step of ranking said possible word choice by the time elapsed since said possible word choice was last used, wherein the word choice temporally closest to its last usage is ranked first.

58. The method as in claim 57 wherein said step (i) of visually representing said possible word choice comprises displaying said possible word choice by the time elapsed since said possible word choice was last used, wherein the word choice temporally closest to its last usage is ranked first.

FIG. 1.

ANALOG-DIGITAL CONVERTOR 4

from microphone 18 (static or dynamic)

from line 20

from telephone 22

amplifier 24

amplifier 24

interface 26

low pass analog filter 28

analog-digital convertor 30

to spectral processor 8

FIGURE 2

SPECTRAL PROCESSER
8

from A-D converter 4

pre-emphasis filter 32

noise model 34

sampler 36

time domain processer 38

MEL scale acoustic band filters 40

frequency domain processer 42

to phoneme labeller 10

FIGURE 3

PHONEME LABELLER
10

from spectral processor 8 ○—

| forward pass labelling | backward pass labelling | generation of machine-readable phonetic symbols |
|:---:|:---:|:---:|
| 44 | 46 | 48 |

—○ to software pre-processor 11

phoneme models 14

training data 16

ASCII look-up tables 15

FIGURE 4

SPEECH RE-SYNTHESIZER 13



FIGURE 5

input$_i$+4
output$_i$
signal$_i$-3
error$_i$-5

FIG. 6a.

input$_i$+4
output$_i$
signal$_i$-3
error$_i$-5

FIG. 6b.

input$_i$+4
output$_i$
signal$_i$-3
error$_i$-5

FIG. 6c.

input$_i$+4
output$_i$
signal$_i$-3
error$_i$-5

FIG. 6d.

$ns := 511$

$i := 0..ns$

$$mean := \frac{\sum\limits_{k=0}^{ns} sample_k}{ns+1} \qquad mean = 3.298 \cdot 10^4 \quad \text{COMPUTED ZERO REFERENCE}$$

### SAMPLE DATA AND MEAN VALUE PLOT



$range := 2^{16}$

$$signal_i := \frac{sample_i - meanv}{range} \qquad \text{ESTIMATED NORMALIZED SIGNAL}$$

## FIG. 7a.

### ABSOLUTE POWER MODEL

$$power1 := \frac{\sum\limits_{k=0}^{ns} |signal_k|}{ns+1}$$

$power1 = 0.142$

### SQUARED POWER MODEL

$$power2 := \sqrt{\frac{\sum\limits_{k=0}^{ns} (signal_k)^2}{ns+1}}$$

$power2 = 0.172$

## FIG. 7b.                    FIG. 7c.

# FIG. 8a.

ESTIMATED SIGNAL PLOT



SIGNAL$_i$ ——
POWERI – –
POWER2 ——

# FIG. 8b.

$$sign(i) := if(signal_i \cdot signal_{i+1} > 0, 0, 1)$$

$$zeros := \frac{\displaystyle\sum_{k=0}^{ns-1} if(sign(k), 0, 1)}{2}$$

zeros = 154.5　　ZERO CROSSING FREQUENCY

SAMPLE DATA AND MEAN VALUE PLOT



FIG. 9a.

$dt := 0.00005$

$j := 0..ns$         $freq_j := \dfrac{j}{dt \cdot ns}$         $win_j := 0.54 - 0.46 \cdot \cos\left(2 \cdot \pi \cdot \dfrac{j}{ns}\right)$

$signal_j := sample_j \cdot win_j$

FIG. 9b.

TAPERED SAMPLE DATA WINDOW



FIG. 9c.

# FIG. 9d.

$$psd := \left[ \left( \overline{|fft(signal)|} \right)^2 \right]$$

$$k := 0.. \frac{ns}{2} \quad lpsd_k := log(psd_k)$$

ESTIMATED POWER SPECTRAL DENSITY

# FIG. 9e.
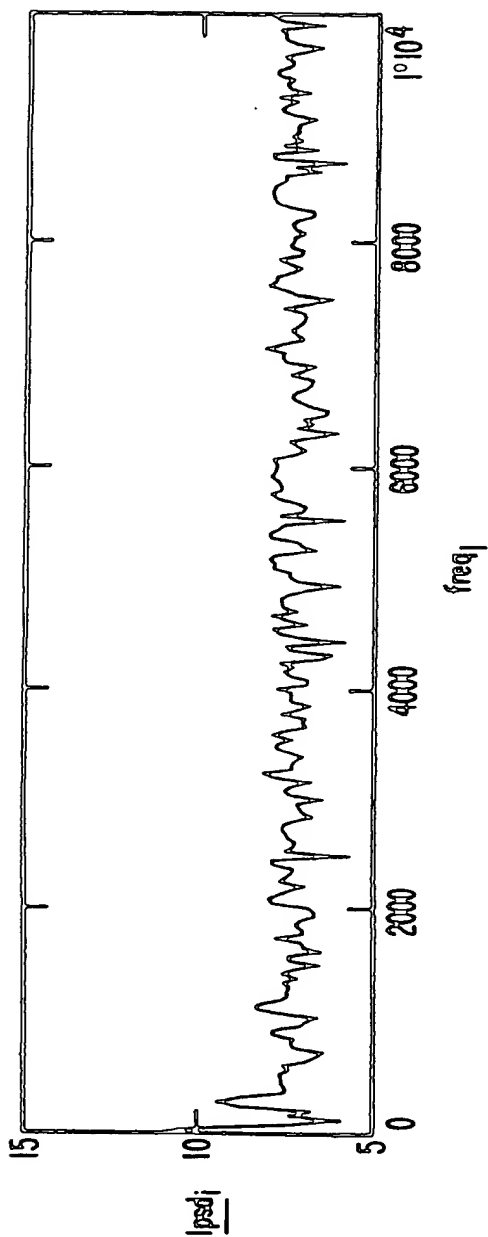
FORMANTS  $fdata_0 = 270$   $fdata_1 = 2.29 \cdot 10^3$   $fdata_2 = 3.07 \cdot 10^3$

$$1 := 0.. \frac{ns}{2}$$

FIG. 9f.

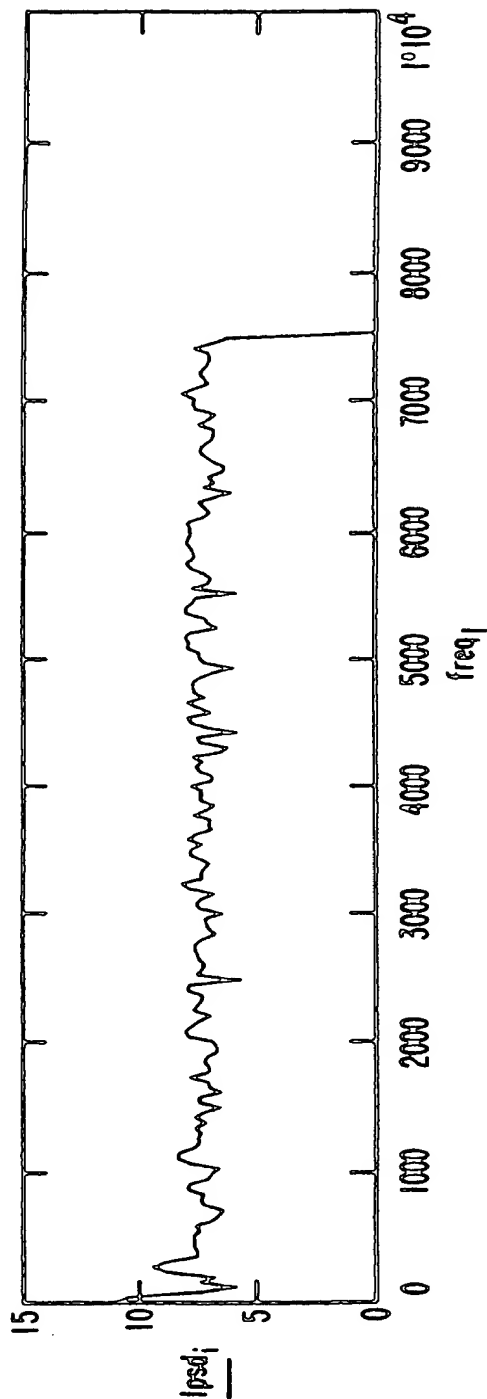$m := 192..ns \quad 1psd_m.0$



FIG. 9g.

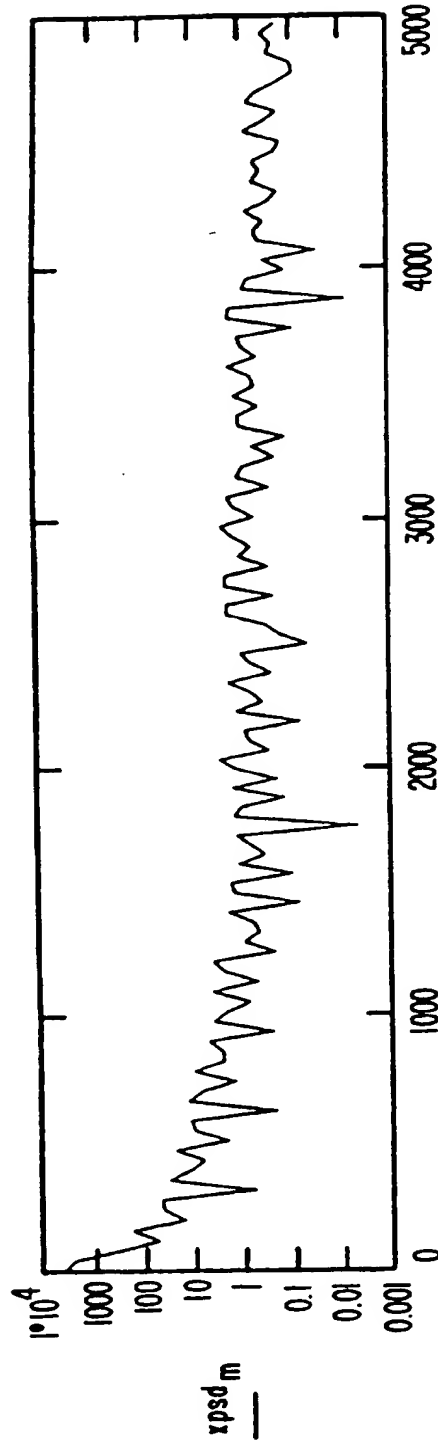$m := 0..127 \quad len := length(lpsd) \quad len=512 \quad ns=511$

$$xpsd := \overline{[(\,|fft(lpsd)|\,)^2]}$$

FIG. 9h.

nFeat :=3   NUMBER OF FEATURES TO COMBINE

nPhon :=5   NUMBER OF PHONEME ALTERNATIVES

nRes :=10   DISCRETE RESOLUTIONS OF MEMBERSHIP FUNCTION

SAMPLE DATA FEATURE VALUES     $f := \begin{pmatrix} 2 \\ 1 \\ 3 \end{pmatrix}$

AVERAGE OF FEATURE VALUES     $fc := mean(f)$     $fc = 2$

PHONEME FEATURE VALUES     $p := \begin{pmatrix} 2 & 2 & 1 & 1 & 3 \\ 1 & 3 & 2 & 3 & 2 \\ 3 & 2 & 4 & 1 & 2 \end{pmatrix}$

AVERAGE OF PHONEME VALUES     $i := 0 .. nPhon-1$

$pc_i := mean(p^{<i>})$     $pc = \begin{bmatrix} 2 \\ 2.333 \\ 2.333 \\ 1.667 \\ 2.333 \end{bmatrix}$

COMPUTE JOINT MEMBERSHIP FUNCTION FOR EACH PHONEME

del :=1   eps :=0.0

$uu_i := wx(fc-del, fc-eps, fc+eps, fc+del, pc_i-del, pc_i-eps, pc_i+eps, pc_i+del, 1, 1)$

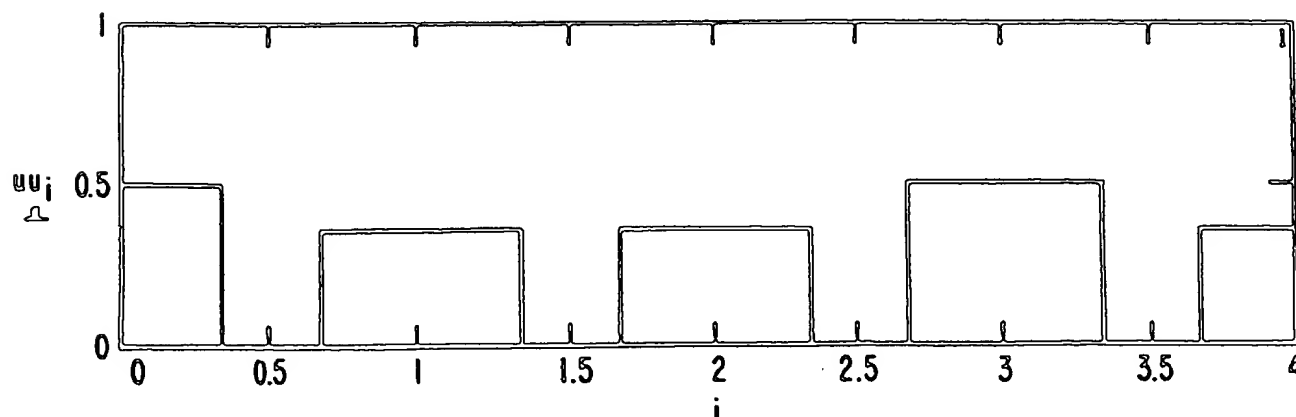POSSIBILITIES FOR EACH PHONEME     $uu = \begin{bmatrix} 0.5 \\ 0.36 \\ 0.36 \\ 0.5 \\ 0.36 \end{bmatrix}$



FIG. 10a.

FIG. 10b.

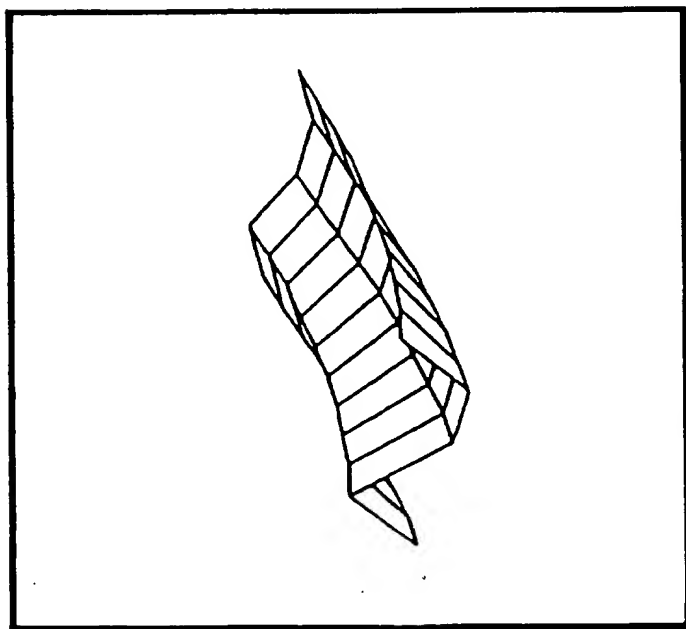COMPUTE VALUES OF PHONEME POSSIBILITY FUNCTIONS FOR DISPLAY

$k := 0..nRes-1$

$$u_{i,k} := wx\left(fc-1, fc, fc, fc+1, pc_i-1, pc_i, pc_i, pc_i+1, 1, 1, \frac{k}{nRes}\right)$$

VALUES OF JOINT MEMBERSHIP FUNCTION

$$u = \begin{bmatrix} 0.2 & 0.217 & 0.236 & 0.257 & 0.281 & 0.308 & 0.338 & 0.372 & 0.41 & 0.452 \\ 0.155 & 0.167 & 0.18 & 0.195 & 0.211 & 0.229 & 0.25 & 0.273 & 0.298 & 0.327 \\ 0.155 & 0.167 & 0.18 & 0.195 & 0.211 & 0.229 & 0.25 & 0.273 & 0.298 & 0.327 \\ 0.2 & 0.217 & 0.236 & 0.257 & 0.281 & 0.308 & 0.338 & 0.372 & 0.41 & 0.452 \\ 0.155 & 0.167 & 0.18 & 0.195 & 0.211 & 0.229 & 0.25 & 0.273 & 0.298 & 0.327 \end{bmatrix}$$

CONTOUR PLOT OF MEMBERSHIP VALUES

SURFACE PLOT OF MEMBERSHIP FUNCTIONS
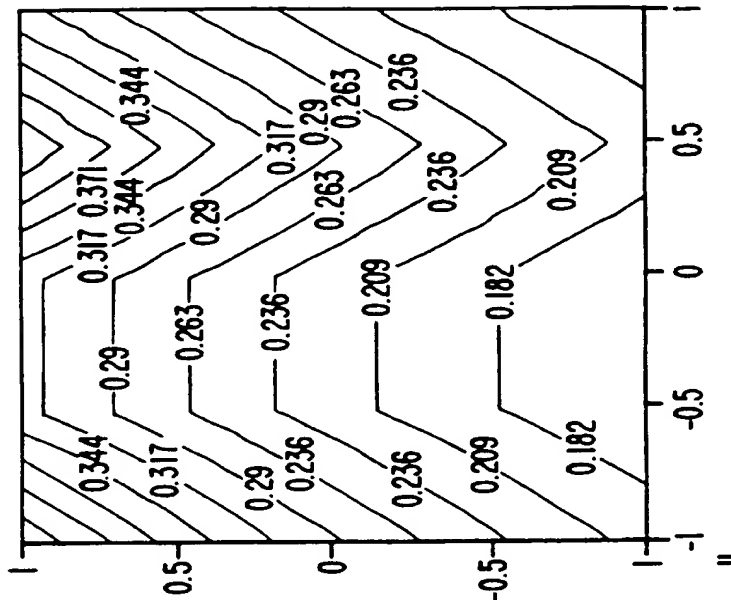
DOCID: <WO___9603741A1_IA>

NUMBER OF PHONEMES            $np := 5$

NUMBER OF SAMPLE DATA BLOCKS    $ns := 9$

POSSIBILITY MATRIX FOR SAMPLE BLOCKS

$$pm := \begin{bmatrix} .4 & .36 & .4 & .14 & .06 & .06 & .05 & .04 & .04 & .04 \\ .11 & .09 & .15 & .35 & .4 & .36 & .31 & .1 & .11 & .11 \\ .04 & .06 & .06 & .11 & .04 & .05 & .15 & .35 & .36 & .4 \\ .3 & .34 & .3 & .09 & .11 & .15 & .11 & .06 & .06 & .06 \\ .09 & .11 & .05 & .16 & .3 & .34 & .29 & .15 & .09 & .09 \\ .06 & .04 & .04 & .14 & .09 & .04 & .09 & .3 & .34 & .3 \end{bmatrix}$$

DETERMINE THE TOP TWO RANKED PHONEME
ASSIGNMENTS OF EACH SAMPLE BLOCK

$i := 0..ns$      $pmax_i := max\left(pm^{<i>}\right)$

$j := 0..np$      $parg_{j,i} := \dfrac{pm_{j,i}}{pmax_i}$      $pass_{j,i} := if\left(parg_{j,i} < 1.0, 0, 1\right)$

$pm_{j,i} := if\left(pass_{j,i} > 0, 0, pm_{j,i}\right)$

$i := 0..ns$      $pmax_i := max\left(pm^{<i>}\right)$

$j := 0..np$      $parg_{j,i} := \dfrac{pm_{j,i}}{pmax_i}$      $pass_{j,i} := if\left[\left(parg_{j,i} < 1.0\right), pass_{j,i}, 2\right]$

THE ASSIGNMENT MATRIX FOR THE TOP TWO
CHOICES IN THE SET OF SAMPLE BLOCKS IS

$$pass = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 2 & 2 & 2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 2 & 2 & 2 \end{bmatrix}$$

## FIG. 11.

## A. CLASSIFICATION OF SUBJECT MATTER
G 10 L 5/06,G 10 L 7/08,G 10 L 9/06,G 10 L 9/18

6

According to International Patent Classification (IPC) or to both national classification and IPC

## B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

G 10 L,H 04 M

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

## C. DOCUMENTS CONSIDERED TO BE RELEVANT

| Category * | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
|---|---|---|
| X | EP, A, 0 248 593 (SPEECH SYSTEMS, INC.) 09 December 1987 (09.12.87), fig. 1; abstract; claims 1-30. -- | 1-10, 44,47-51 |
| X,P | EP, A, 0 645 757 (XEROX CORP.) 29 March 1995 (29.03.95), fig. 1; abstract; claims 1-12. -- | 1-58 |
| A | US, A, 5 289 523 (VASILE et al.) 22 February 1994 (22.02.94), fig. 3; abstract; claims 1-21. ---- | 1-58 |

☐ Further documents are listed in the continuation of box C.

☐ Patent family members are listed in annex.

* Special categories of cited documents :

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier document but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.

"&" document member of the same patent family

| Date of the actual completion of the international search | Date of mailing of the international search report |
|---|---|
| 12 October 1995 | 1 0. 11. 95 |

| Name and mailing address of the ISA | Authorized officer |
|---|---|
| European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Tx. 31 651 epo nl, Fax (+31-70) 340-3016 | BERGER e.h. |

Form PCT/ISA/210 (second sheet) (July 1992)

# ANHANG

## ANNEX

## ANNEXE

zum internationalen Recherchen-
bericht über die internationale
Patentanmeldung Nr.

to the International Search
Report to the International Patent
Application No.

au rapport de recherche inter-
national relatif à la demande de brevet
international n°

## PCT/US 95/09130 SAE 114738

In diesem Anhang sind die Mitglieder
der Patentfamilien der im obenge-
nannten internationalen Recherchenbericht
angeführten Patentdokumente angegeben.
Diese Angaben dienen nur zur Unter-
richtung und erfolgen ohne Gewähr.

This Annex lists the patent family
members relating to the patent documents
cited in the above-mentioned inter-
national search report. The Office is
in no way liable for these particulars
which are given merely for the purpose
of information.

La présénte annexe indique les
membres de la famille de brevets
relatifs aux documents de brevets cités
dans le rapport de recherche inter-
national visée ci-dessus. Les renseigne-
ments fournis sont donnés à titre indica-
tif et n'engagent pas la responsibilité
de l'Office.

| Im Recherchenbericht angeführtes Patentdokument Patent document cited in search report Document de brevet cité dans le rapport de recherche | Datum der Veröffentlichung Publication date Date de publication | Mitglied(er) der Patentfamilie Patent family member(s) Membre(s) de la famille de brevets | Datum der Veröffentlichung Publication date Date de publication |
|---|---|---|---|
| EP A1 248590 | 09-12-87 | JP A2 63066600 | 25-03-88 |
| | | US A 5054085 | 01-10-91 |
| | | AT E 45831 | 15-09-89 |
| | | AU A1 28312/84 | 21-11-85 |
| | | AU B2 590000 | 22-12-89 |
| | | CA A1 1216944 | 20-01-87 |
| | | DE C0 3479543 | 28-09-89 |
| | | EP A1 126420 | 28-11-84 |
| | | EP B1 126420 | 23-08-89 |
| | | JP A2 59221198 | 12-12-84 |
| | | US A 4718096 | 05-01-88 |
| EP A1 645757 | 29-03-95 | JP A2 7175497 | 14-07-95 |
| US A 5289523 | 22-02-94 | keine - none - rien | |